Self-Supervised Learning in Vision

Slides from Ishan Misra, Naiyan Wang & many others.

Lecture 6

Rob Fergus

Success story of supervision

ImageNet Challenge

Classification Results (CLS)





Success story of supervision: Pre-training Features from networks pre-trained on ImageNet can be used for a variety of different downstream tasks



Images from ImageNet (Pre-train)

Learn a representation

facebook Artificial Intelligence Research ConvNet



Success story of supervision: Recipe for good solutions

- Pre-train on a large supervised dataset.
- Collect a dataset of "supervised" images
- Train a ConvNet



Can we get labels for all data?

facebook Artificial Intelligence Research





Bounding Boxes

Stats from Pawan Kumar at Oxford







Image Level



7



Artificial Intelligence Research

https://www.forbes.com/sites/bernardmarr/2018/05/21/how-much-data-do-we-create-every-day-the-mind-blowing-stats-everyone-should-read/





Bounding Boxes

Image Level

Internet Photos





Bounding Boxes

ImageNet (14 million images) needed 22 human years to label

Image Level

Internet Photos



Can we get labels for all

data? about complex concepts?

- Video?
- Labelling cannot scale to the size of the data we generate

11

Rare concepts?



Objects in Vision Dataset (LabelMe)

		-					_	-					-					•			=									
		*	ć,	4	÷	÷.	1	e,	÷	ŝ.	Ċ.	ç,	ŝ	γ.		۰.	1		č,	ì	'n,		۰.	-		ę.	ŝ		-	
		,	ć	į,	÷	1	1	÷	÷	÷	ł	ç	÷	1		÷		4	7	;	ł		٦	4		ŝ	ł			
1		• *	e'	\mathbf{r}^{\prime}				e'	÷	-	1	÷	•	٦,				•	Ξ.	1	1		٠,	•	•	2	ł			
-			. 1	. 1				. :		• •	. 1	• 1	• •	1			1	• :	1		1		• •			1	-			
											1	1				ŝ	1			1	1	1			2	1			1	
		÷.,	1.	2	-			÷.,			i.	ŝ-	-			5	5	5	7	1	ľ		۰.	-		ŝ-				
n		-	-	-	5	-	_	-				-				-	-	-	5	Ξ	:		-	-	-	-	Ī		-	
					÷							1								-	:					÷	i			
		٠	÷	I.	÷	-	2	;			1		÷	2		1	ı	:	:	÷	ł		۱	٠	•	1	i	2	•	
		,	1	,	÷	i		•				÷	:				÷		;	÷	:		,	1		ł	į			
			÷	. 1	- 1		1	1			2	- 2	• 3	۰.		۰.	•		1	1	1		•	1		2				
1											1				1		1			1	ľ			-						
		-	1	۰,	Ξ.	1	1		ĩ.	ē,	2	5	-	÷.		1			1	5	e'		1	1	1	į,				_
,			1	1	- 2			Ċ,	÷	÷	2	9	5	÷.		1		1	7	÷	÷			1		2	į			_
		• *	d,	d.			C	Ċ,		÷	Ŷ	Ŷ		÷,		2	1.0	12	7	1	Y		12	• 14		2	ŝ			
		- 5	J.					- 1			- 6	÷	-			1	- 2	-7		÷.	Ē		- 1			÷	- 1			
3			_					•					:							÷	:			1		÷	į		•	
			ŀ									÷	:							÷	:					÷	i			
	•	D	e ²	\mathcal{L}^{2}				·		•	1	- 2	• 2	۰.		۰.	٠.	٠.	1		1		•	•	•	9	ł			
		۰,	5	\mathbf{x}		Ĩ,	ć	ī,			2	2	2	2		1			7	÷	ŕ					2	5			
		•••	5	\mathcal{L}		0		2	-		11	23	1	÷.,	-	4	÷**	÷.,		1	ċ,	-	1	1		23	5			
	L				ł			t			t	ł	F				ł	ł		ł						ł			t	L

10% of the classes account for 93% of the data

Slide credit: Rob Fergus





Different Domains?



facebook Artificial Intelligence Research

ImageNet pre-training may not work

13

Arguments for Unsupervised Learning

- Want to be able to exploit unlabeled data
 - Vast amount of it often available
 - Essentially free
- Good regularizer for supervised learning
 - Helps generalization
 - Transfer learning
 - Zero/one/few shot learning

Unsupervised Learning

- Biological argument [from G. Hinton]:
 - Our brains have 10^15 connections
 - We live for 10^9 secs
 - Need 10^6 bits/sec

 - Insufficient information from occasional high level label Only source with enough information is input itself
- Challenging problem: big focus on many DL groups

Historical Note

- Deep Learning revival started in ~2006 Hinton & Salakhudinov Science paper on RBMs
- Unsupervised Deep Learning was focus from 2006-2012
- In ~2012 great results in vision, speech with supervised methods appeared

 - Initially less interest in unsupervised learning • By focus once more on unsupervised learning

Overview of Unsupervised Perspectives

- Given just data {X}
 - Unlike supervised learning there are no **provided** labels {Y}
- 1. Density modeling, i.e. build model of p(X)
 - Enables sampling of new data
 - Evaluate probability of a data point
 - Can be conditional model, e.g. p(X_t | X_{t-1},...)
 - Requires (deep) generative architectures [HARD]
 - Generating pixels not necessarily optimal for learning representations for downstream tasks

1. Density Modeling

- Have access to $x \sim p_{data}(x)$ through training set
- Want to learn a model $x \sim p_{model}(x)$
- Want p_{model} to be similar to p_{data} :

Samples from true data distribution have high likelihood under p_{model}



Samples drawn from p_{model} reflect structure of p_{data}



2. "Self supervised" learning

- Also unsupervised but...
- Find supervision signal y within the input data
- $y: \mathcal{X} \to \mathcal{Y}$ $x \mapsto y(x)$
- architectures

$$\min_{\theta} \frac{1}{n} \sum_{i=1}^{n}$$

- Pre-training of representation for subsequent task
- Typically involves some insight into domain to pick y

• This signal is then used as a target in *discriminative model*:

Allows the use of standard supervised learning losses and

 $\ell(f_{\theta}(x_i), y(x_i))$

Class Project

- during office hours.
- Simple option is to reimplement existing approach/paper.
 - Cannot just copy existing repo, but can use it for debugging.
- Alternatively, you can extend existing approach/paper.
 - Can use existing repo, but extension must be significant
- More ambitious (more marks): Implement new idea from scratch
- Can use existing models/code but must have significant novelty Project abstracts due Oct 15th (email to me).

 - Email paragraph summarizing your project plan to <u>fergus@cs.nyu.edu</u>. • Please be sure to include the names of people on your project.

Everyone must come and discuss proposed ideas with me at some point

Class Project details

Team of 2-3 people [1 person not allowed; 3 is hard max] PyTorch preferred

Project video:

- 2 minute clip explaining your project
- Voice over slides

Project Report

- Format: 4-8 page conference paper style report on your project (please don't waffle)
- Intro (with refs to related work) [~1 page]
- Method (be sure to cite any code/pre-trained models) [~2 pages]

- Experiments (must have plots/results figures; also should have baselines; ideally some kind of ablation experiments too) [~2-4 pages]

- Discuss (brief) [~0.5 pages
- See examples: http://openaccess.thecvf.com/CVPR2018.py
- Zip of source code or link to Github (please ensure you give access to robfergus)

Class Project details (2)

- . Deadline is Thursday Dec 17th midnight [Hard deadline]
- Feel free to turn in earlier
- Will try to grade them and compute final grades by Christmas
- Grading (49% of total grade for class)
- Novelty / Technical difficulty of problem [15%]
- Quality of Results [15%]
- Quality of implementation [5%]
- Quality of writeup & video presentation [14%]
- How many people in your group

Class Project General Advice

- Please make sure you have *something* working, even if you don't achieve overall goal

- Even a small part of an ambitious project can be OK
- So please have a safe plan B option in mind
- Expect all projects to train something, i.e. must use b-prop at some point - Just evaluating existing models is NOT OK.

moment.

- Cluster gets busy at end of semester -- please don't leave it all to last



Self-supervised learning in computer vision

facebook Artificial Intelligence Research Ishan Misra

With slides from Andrew Zisserman, Carl Doersch





What is "self" supervision?

• Obtain labels from the data itself by using a "semi-automatic" process





What is "self" supervision?

• Obtain labels from the data itself by using a "semi-automatic" process



implausible label





Unsupervised

- limited power

Self-Supervised - derives label from a co-occuring input to another modality







• Fill in the blanks

Softmax classifier



Hidden layer

Projection layer



facebook Artificial Intelligence Research word2vec - Mikolov et al. Image by Julian Gilyadov 27

Success of self-supervised learning in NLP

- Fill in the blanks is a powerful signal to learn representations
- Sentence/Word representations: BERT Devlin et al., 2018



Why self supervision?

- Helps us learn using observations and interactions Does not require exhaustive annotation of concepts • Leverage multiple modalities or structure in the domain





images

- Doersch et al. (2015)
- Zhang et al. (2016)
- Zhang et al. (2017)
- Noroozi et al. (2016)
- Pathak et al. (2016)

richer data





images

videos richer data

- Doersch et al. (2015)
- Zhang et al. (2016)
- Zhang et al. (2017)
- Noroozi et al. (2016)
- Pathak et al. (2016)

- Wang et al. (2015)
- Misra et al. (2016)
- Pathak et al. (2017)







images

videos sound & depth richer data

- Doersch et al. (2015)
- Zhang et al. (2016)
- Zhang et al. (2017)
- Noroozi et al. (2016)
- Pathak et al. (2016)

- Wang et al. (2015)
- Misra et al. (2016)
- Pathak et al. (2017)



- Owens et al. (2016)
- Zhang et al. (2017)
- Bansal et al. (2016)





images

sound & depth videos richer data

- Doersch et al. (2015)
- Zhang et al. (2016)
- Zhang et al. (2017)
- Norouzi et al. (2016)
- Pathak et al. (2016)

- Wang et al. (2015)
- Misra et al. (2016) • Pathak et al. (2017)





actions

- Owens et al. (2016)
- Zhang et al. (2017)
- Bansal et al. (2016)
- Agarwal et al. (2015)
- Jayaraman et al. (2015)
- Pinto et al. (2016)
- Agarwal et al. (2016)
- Pinto et al. (2017)
- Pinto et al. (2016)

Self-supervision in computer vision

- Using images
- Using video
- Using video and sound



From Images: Relative position of patches



Sample Second Patch

Unsupervised visual representation learning by context prediction, Carl Doersch, Abhinav Gupta, Alexei A. Efros, ICCV 2015



From Images: Relative position of patches





facebook Artificial Intelligence Research

Unsupervised visual representation learning by context prediction, Carl Doersch, Abhinav Gupta, Alexei A. Efros, ICCV 2015




Relative Position: Nearest Neighbors in features



Relative-positioning Input



Random Initialization

ImageNet AlexNet

Unsupervised visual representation learning by context prediction, Carl Doersch, Abhinav Gupta, Alexei A. Efros, ICCV 2015





Unsupervised Learning of Visual Representations by Solving Jigsaw Puzzles, Noorozi et al. (2016)

What do we learn when we solve a Jigsaw puzzle?

[Noorozi et al. (2016)]



Hash Set

7

index	table	Re
		hc
64	9,4,6,8,3,2,5,1,7	

eorder patches ccording to the selected ash table



[Noorozi et al. (2016)]



[Noorozi et al. (2016)]

Visualization of filters



Visualization of the top 16 activations for 6 selected channels of the convolutional layers

Noroozi, M., & Favaro, P. Unsupervised learning of visual representations by solving jigsaw puzzles. In ECCV 2016.



Feature Learning by Inpainting

[Context Encoders: Feature Learning by Inpainting, Pathak et al. (2016)]





[Pathak et al. (2016)]

Context Encoders





- Encoder can be substituted with any network architecture like AlexNet etc.
- Decoder is a set of UpConv/deconv/frac-stridedconv layers

[Pathak et al. (2016)]

Combined L2 + GAN loss



Input Image L2 Loss



Adversarial Loss Joint Loss

[Pathak et al. (2016)]

Image colorization



Colorful Image Colorization Richard Zhang, Phillip Isola, Alexei (Alyosha) Efros

http://richzhang.github.io/colorization/

slides from Zhang



$\begin{array}{lll} \textbf{Gravscale image: } L \text{ channel} & \textbf{Color information: } ab \text{ channels} \\ \mathbf{X} \in \mathbb{R}^{H \times W \times 1} \end{array}$

 \mathcal{F}





$$\rightarrow$$
 ab

Semant

$\begin{array}{c} \text{Gravscale image} \\ \mathbf{X} \in \mathbb{R}^{H \times W} \end{array} \begin{array}{c} \text{Seman} \\ \text{level a} \end{array}$





From Images: Predicting Rotations

Which image has the correct rotation?







facebook Artificial Intelligence Research



48

[Dosovitskiy et al. ICLR 2014]



- 1 class = single image + its transformations
- Learn to classify each "class"
- Domain knowledge about appropriate transformations
- does not scale

Many different self-supervision tasks, how to evaluate?

facebook Artificial Intelligence Research

Context auto encoders - Pathak et al., 2016



Self-supervised pre-training



Pre-train data

facebook Artificial Intelligence Research





Learn a representation



51

Fine-tune on end task (Image Classification)





Tests representation as well as how good the pre-training initialization is



Fine-tune on end task

Initialization (ResNet101)

ImageNet Supervised

Relative Position

Colorization

facebook Artificial Intelligence Research

End task ImageNet top-5 accuracy VOC07 Detection mAP 85.1 74.2 66.8 59.2 62.5 65.5

• Multi-task self-supervised visual learning, C Doersch, A Zisserman, ICCV 2017



53

Are they complementary?

Initialization (ResNet101)

ImageNet Supervised

Relative Position

Colorization

Relative Position + Colorization (Multi-task)

facebook Artificial Intelligence Research



End lask

ImageNet top-5 accuracy VOC07 Detection mAP

85.1	74.2
59.2	66.8
62.5	65.5
66.6	68.8

• Multi-task self-supervised visual learning, C Doersch, A Zisserman, ICCV 2017



54

Train linear classifiers on "fixed" features



Tests how good the representation is (linearly separable)



Self-supervision in computer vision

- Using images
- Using video
- Using video and sound



Video

- Video is a "sequence" of frames
- How to get "self-supervision"?

- Predict order of frames
- Fill in the blanks
- Track objects and predict their position



"Sequence" of data



Video

- Slow feature
 - Neighborhood frames should have similar features

 $\mathcal{U}_2 = \{ \langle (j,k), p_{jk} \rangle : x_j, x_k \in \mathcal{U} \}$ $R_2(\boldsymbol{\theta}, \mathcal{U}) = \sum D_{\delta}(\mathbf{z}_{\boldsymbol{\theta}}(\boldsymbol{x}_j), \mathbf{z}_{\boldsymbol{\theta}}(\boldsymbol{x}_k), p_{jk})$ $(j,k) \in \mathcal{U}_2$ = $\sum p_{jk} d(\mathbf{z}_{\theta j}, \mathbf{z}_{\theta k}) + \overline{p_{jk}} \max(\delta - d(\mathbf{z}_{\theta j}, \mathbf{z}_{\theta k}), 0),$ $(j,k) \in \mathcal{U}_2$

[A Survey to Self-Supervised Learning, Naiyan Wang]

$$\mathcal{I} \text{ and } p_{jk} = \mathbb{1}(0 \le j - k \le T) \},$$

Mobahi, H., Collobert, R., & Weston, J. Deep learning from temporal coherence in video. In ICML 2009.

Wiskott, L., & Sejnowski, T. J. (2002). Slow feature analysis: Unsupervised learning of invariances. *Neural computation*, 14(4), 715-770.

Video

- Slow and steady feature
 - Not only similar, but also smooth
 - Extend to triplet setting (Not triplet loss!)



 $D_{\delta}(\mathbf{z}_{\boldsymbol{\theta}l} - \mathbf{z}_{\boldsymbol{\theta}m}, \mathbf{z}_{\boldsymbol{\theta}m} - \mathbf{z}_{\boldsymbol{\theta}n}, p_{lmn}),$ $R_3(\boldsymbol{\theta}, \mathcal{U}) =$ $(l,m,n) \in \mathcal{U}_3$

[A Survey to Self-Supervised Learning, Naiyan Wang]

 $\mathcal{U}_3 = \{ \langle (l, m, n), p_{lmn} \rangle : \boldsymbol{x}_l, \boldsymbol{x}_m, \boldsymbol{x}_n \in \mathcal{U} \text{ and } p_{lmn} = \mathbb{1} (0 \leq m - l = n - m \leq T) \}.$

Jayaraman, D., & Grauman, K. Slow and steady feature analysis: higher order temporal coherence in video. In CVPR 2016.

Temporally Correct order



Temporally Incorrect order





lmages





Given a start and an end, can this point lie in between?

facebook Artificial Intelligence Research





Input Tuple





Nearest Neighbors of Query Frame (fc7 features)

Query

ImageNet









Shuffle & Learn





Random





63

Fine-tune on Human Keypoint Estimation



64

Fine-tune on Human Keypoint Estimation

Initialization (AlexNet)

ImageNet Supervised

Shuffle and Learn

facebook Artificial Intelligence Research







From video: encoding more structure

	Predicted odd element
fc6 conv5	
Video-clip Encoder	Vide
Correct order	×





Unsupervised Learning of Visual Representations using Videos, Wang & Gupta 2015

Idea: Object Tracking in Videos





Approach





(a) Unsupervised Tracking in Videos



(b) Siamese-triplet Network

(c) Ranking Objective

- Use object tracking in videos
 - Classify if patches belong to the same track or not



Patch Mining In Videos

- Track 8M patches in 100K videos from YouTube.
- Use off-the-shelf tracking algorithms with no learning.





VOC 2007 Detection Performance (pretraining for R-CNN)



Precision Average %



Object Movement

- The world is rigid, or at least piecewise rigid
 - Motion provide evidence of how pixels move together
 - The pixels move together are likely to form an object



[A Survey to Self-Supervised Learning, Naiyan Wang]

3. Train ConvNet

Pathak, D., Girshick, R., Dollár, P., Darrell, T., & Hariharan, B. Learning Features by Watching Objects Move. In CVPR 2017.



Self-supervision in computer vision

- Using images
- Using video
- Using video and sound


Audio-Visual co-supervision

Train a network to predict if image and audio clip correspond



Correspond?

facebook Artificial Intelligence Research





Audio-Visual co-supervision

positive







Audio-Visual co-supervision



facebook Artificial Intelligence Research





Audio-Visual co-supervision What can be learnt?



facebook Artificial Intelligence Researc

- Good representations Visual features
- Audio features

- Intra- and cross-modal retrieval
- Aligned audio and visual embeddings

- "What is making the sound?"
- Learn to localize objects that sound



Audio-Visual co-supervision What would make this sound?



Note, no video (motion) information is used

facebook Artificial Intelligence Research



Visual + Audio

[Ambient Sound Provides Supervision for Visual Learning, [Owens et al. (2016)]







Image







- Ego-motion
 - "We move in order to see and we see in order to move" J.J Gibson
 - Ego-motion data is easy to collect
 - axises. (Visual Base-CNN Stream-1



[A Survey to Self-Supervised Learning, Naiyan Wang]

Siamese CNN to predict camera translation & Rotation along 3-

Agrawal, P., Carreira, J., & Malik, J. Learning to see by moving. In ICCV 2015

- Ego-motion
 - Learning features that are equivariant to ego-motion



[A Survey to Self-Supervised Learning, Naiyan Wang]

Jayaraman, D., & Grauman, K. Learning image representations tied to ego-motion. In *ICCV 2015*



- Ego-motion

 - Siamese networks with contrastive loss • M g is the transformation matrix specified by the external sensors

$$(oldsymbol{ heta}^*, \mathcal{M}^*) = rgmin_{oldsymbol{ heta}, \mathcal{M}} \sum_{g, i, j} d_g$$
 $d_g(oldsymbol{a}, oldsymbol{b}, c) = \mathbbm{1}(c = g)d(oldsymbol{a}, b)$

[A Survey to Self-Supervised Learning, Naiyan Wang]

 $l_g \left(M_g \mathbf{z}_{\boldsymbol{\theta}}(\boldsymbol{x}_i), \mathbf{z}_{\boldsymbol{\theta}}(\boldsymbol{x}_j), p_{ij} \right),$

 $(a, b) + \mathbb{1}(c \neq g) \max(\delta - d(a, b), 0),$

Jayaraman, D., & Grauman, K. Learning image representations tied to ego-motion. In ICCV 2015



- Acoustics -> RGB
 - Similar events should have similar sound.
 - Naturally cluster the videos.



[A Survey to Self-Supervised Learning, Naiyan Wang]

Owens, A., Wu, J., McDermott, J. H., Freeman, W. T., & Torralba, A. Ambient sound provides supervision for visual learning. In ECCV 2016

- Features for grasping
 - angle





Query Kinect image



Verify whether we could grasp the center of a patch at a given

Pinto, L., & Gupta, A. Supersizing self-supervision: Learning to grasp from 50k tries and 700 robot hours. In ICRA 2016



Evaluation

Evaluate on general high-level vision tasks (classification, detection) Be caution of different settings!

	Full train set				150 image set								
Method	All	>c1	>c2	>c3	>c4	>c5	All	>c1	>c2	>c3	>c4	>c5	#wins
Supervised													
Imagenet	56.5	57.0	57.1	57.1	55.6	52.5	17.7	19.1	19.7	20.3	20.9	19.6	NA
Sup. Masks (Ours)	51.7	51.8	52.7	52.2	52.0	47.5	13.6	13.8	15.5	17.6	18.1	15.1	NA
Unsupervised													
Jigsaw [‡] [30]	49.0	50.0	48.9	47.7	45.8	37.1	5.9	8.7	8.8	10.1	9.9	7.9	NA
Kmeans [23]	42.8	42.2	40.3	37.1	32.4	26.0	4.1	4.9	5.0	4.5	4.2	4.0	0
Egomotion [2]	37.4	36.9	34.4	28.9	24.1	17.1	_	_	_	_	_	_	0
Inpainting [35]	39.1	36.4	34.1	29.4	24.8	13.4	_	_	_	_	_	_	0
Tracking-gray [46]	43.5	44.6	44.6	44.2	41.5	35.7	3.7	5.7	7.4	9.0	9.4	9.0	0
Sounds [33]	42.9	42.3	40.6	37.1	32.0	26.5	5.4	5.1	5.0	4.8	4.0	3.5	0
BiGAN [10]	44.9	44.6	44.7	42.4	38.4	29.4	4.9	6.1	7.3	7.6	7.1	4.6	0
Colorization [51]	44.5	44.9	44.7	44.4	42.6	38.0	6.1	7.9	8.6	10.6	10.7	9.9	0
Split-Brain Auto [52]	43.8	45.6	45.6	46.1	44.1	37.6	3.5	7.9	9.6	10.2	11.0	10.0	0
Context [8]	49.9	48.8	44.4	44.3	42.1	33.2	6.7	10.2	9.2	9.5	9.4	8.7	3
Context-videos [†] [8]	47.8	47.9	46.6	47.2	44.3	33.4	6.6	9.2	10.7	12.2	11.2	9.0	1
Motion Masks (Ours)	48.6	48.2	48.3	47.0	45.8	40.3	10.2	10.2	11.7	12.5	13.3	11.0	9

[A Survey to Self-Supervised Learning, Naiyan Wang]

Pathak, D., Girshick, R., Dollár, P., Darrell, T., & Hariharan, B. Learning Features by Watching Objects Move. In CVPR 2017.



Main issue with all these methods

- All these models rely on expert knowledge
- Need to define y(x) for each new domain
- Not clear how to select a y(x) the features

Not clear how to select a y(x) that is a good target to learn all-purpose

Unsupervised Learning by Predicting Noise

Piotr Bojanowski, Armand Joulin ICML 2017

Unsupervised Learning by Predicting Noise [Bojanowski & Joulin, ICML 2017]



- Inspired by Dosovitskiy et al.
- Learn mapping from images to a sphere
- Fix targets on sphere
- Simultaneously:
 - Learn the mapping
 - Optimize the assignment between images and targets

Deep Discriminative Clustering

- We are given a set of n images
- We want to learn a visual features f without using labels



min min

• We use the L2 loss

$$\{x_1,\ldots,x_n\}$$

$$\min_{y_i} \ell(f_\theta(x_i), y_i)$$

$$\frac{1}{n} \|f_{\theta}(X) - Y\|_F^2$$

Label Collapse Problem

- Optimization over Y would lead to a collapse
- Repulsive costs are tricky to use

Can impose constraints on Y but hard to optimize

Fixing the Target Representation

- Instead, we fix the target representation
- Allow a reassignment between targets and images

$$Y = PC \qquad \qquad \mathcal{P} = \{I$$

• Targets C are uniformly sampled on the sphere

$$\min_{\theta} \min_{P \in \mathcal{P}} \frac{1}{2n} \| f_{\theta}(X) - PC \|_{F}^{2}$$
• Final objective function

 $P \in \{0,1\}^{n \times k} \mid P\mathbf{1} = \mathbf{1}, P^{\top}\mathbf{1} = \mathbf{1}\}$



Optimization

- We minimize our cost function in an on-line fashion
- We use the following algorithm:

Require: T batches of images, $\lambda_0 > 0$ for $t = \{1, ..., T\}$ do Compute $f_{\theta}(X_b)$ Compute $\nabla_{\theta} L(\theta)$ using P^* Update $\theta \leftarrow \theta - \lambda_t \nabla_{\theta} L(\theta)$ end for

Obtain batch b and representations r

Compute P^* by minimizing w.r.t. P

Optimizing the Permutation Matrix

- At theta fixed, the permutation is obtained by solving
 - $\max_{P \in \mathcal{P}} \operatorname{Tr} ($
- Which is a linear program on the set of permutation matrices • We can use the Hungarian algorithm

$$\left(PCf_{\theta}(X)^{\top}\right).$$

 $\mathcal{O}(nb^2)$

Experimental Setup



- AlexNet architecture
- set
- i.e. PASCAL VOC Classification / Detection

Learn unsupervised features on ImageNet training

Retrain a classifier on top for a target transfer task, [Bojanowski & Joulin, ICML 2017]

Baselines

- Self supervised models
 - Wang & Gupta Temporal coherence in videos
 - Doersch et al. Predict context patches
 - Zhang et al. Predict color
 - Norouzi & Favaro Solve jigsaw puzzles
- Unsupervised model
 - GAN
 - Auto-encoder
 - BI-GAN (Donahue et al.)

Pascal VOC - results

	Classification		Detection	
Trained layers	fc6-8	all	all	
ImageNet labels	78.9	79.9	56.8	
Agrawal et al.	31.0	54.2	43.9	
Pathak et al.	34.6	56.5	44.5	
Wang & Gupta	55.6	63.1	47.4	
Doersch et al.	55.1	65.3	51.1	
Zhang et al.	61.5	65.6	46.9	
Autoencoder	16.0	53.8	41.9	
GAN	40.5	56.4	-	
BiGAN	52.3	60.1	46.9	
NAT	56.7	65.3	49.4	



 Poor performance of AE / GAN

Nearest Neighbor Queries









Bojanowski & Joulin Summary

- Simple unsupervised approach
- No domain expert knowledge
- Scales to very large datasets
- Close to supervised pipeline
- SOTA performance (at the time) amongst unsupervised methods

Contrastive Predictive Coding

Aaron van den Oord, Yazhe Li, Oriol Vinyals

Google DeepMind

https://arxiv.org/pdf/1807.03748.pdf

Model Overview



CPC Principle

signals x and c defined as

$$I(x;c) = \sum_{x,c} p(x,c) \log \frac{p(x|c)}{p(x)}$$

 By maximizing the mutual information between the encoded common.

• Encode the target x (future) and context c (present) into a compact distributed vector representations (via non-linear learned mappings) in a way that maximally preserves the mutual information of the original

representations (which is bounded by the MI between the input signals), we extract the underlying latent variables the inputs have in

CPC Model

- Do NOT predict future observations x_{t+k} directly with a generative model $p_k(x_{t+k} | c_t)$
- Instead we model a <u>density ratio</u> which preserves the mutual information between x_{t+k} and c_t (prev eqn):

$$f_k(x_{t+k}, c_t) \propto \frac{p(x_{t+k}|c_t)}{p(x_{t+k})}$$

Where: $f_k(x_{t+k}, c_t) =$

$$= \exp\left(z_{t+k}^T W_k c_t\right),$$

Noise Contrastive Estimation Loss

$$\mathcal{L}_{N} = -\mathbb{E}_{X} \left[\log \frac{f_{k}(x_{t+k}, c_{t})}{\sum_{x_{j} \in X} f_{k}(x_{j}, c_{t})} \right]$$

ratio in prev slide.

- - E.g. same/different sequence? Narrow window?

• Given a set $X = \{x_1, \ldots, x_N\}$ of N random samples containing one positive sample from $p(x_{t+k} \mid c_t)$ and N – 1 negative samples from the 'proposal' distribution $p(x_{t+k})$, we optimize

Optimizing this loss will result in $f_k(x_{t+k}, c_t)$ estimating the density

Not always clear in expts where the N random samples come from

C.F. Vision SSL approaches

[Wang & Gupta 2015]



(a) Unsupervised Tracking in Videos



Tracked Negative Query (First Frame) (Last Frame) (Random) (b) Siamese-triplet Network



D: Distance in deep feature space

(c) Ranking Objective

Unsupervised Visual Representation Learning by Context Prediction [Doersch et al. ICCV 2015]



CPC applied to Audio

Method	ACC	Mot
Phone classification Random initialization MFCC features CPC Supervised	27.6 39.7 64.6 74.6	#step 2 ste 4 ste 8 ste
Speaker classification Random initialization MFCC features CPC Supervised	1.87 17.6 97.4 98.5	- 12 st 16 st Nega Mixe Same Mixe

Table 1: LibriSpeech phone and speaker classification results. For phone classification there are 41 possible classes and for speaker classification 251. All models used the same architecture and the same audio input sizes.

Method	ACC
#steps predicted	
2 steps	28.5
4 steps	57.6
8 steps	63.6
12 steps	64.6
16 steps	63.8
Negative samples from	
Mixed speaker	64.6
Same speaker	65.5
Mixed speaker (excl.)	57.3
Same speaker (excl.)	64.6
Current sequence only	65.2



Figure 2: t-SNE visualization of audio (speech) representations for a subset of 10 speakers (out of 251). Every color represents a different speaker.



Figure 3: Average accuracy of predicting the positive sample in the contrastive loss for 1 to 20 latent steps in the future of a speech waveform. The model predicts up to 200ms in the future as every step consists of 10ms of audio.

 Table 2: LibriSpeech phone classifica tion ablation experiments. More details can be found in Section 3.1.

CPC applied to images



	Method	Top-1 ACC
	Using AlexNet conv5	
	Video [27]	29.8
	Relative Position [11]	30.4
	BiGan [34]	34.8
	Colorization [10]	35.2
	Jigsaw [28] *	38.1
. •	Using ResNet-V2	
tions	Motion Segmentation [35]	27.6
	Exemplar [35]	31.5
	Relative Position [35]	36.2
	Colorization [35]	39.6
	CPC	48.7

Table 3: ImageNet top-1 unsupervised classification results. *Jigsaw is not directly comparable to the other AlexNet results because of architectural differences.
Google Research

SimCLR: A Simple Framework for Contrastive Learning of Visual Representations

- Ting Chen
- Simon
- Kornblith
- Mohammad Norouzi
- **Geoffrey Hinton**

Google Research, Brain Team



A simple idea: maximizing the agreement of representations under data transformation, using a contrastive loss in the latent/feature space.



Google Rese

Figure 2. A framework for contrastive representation learning. Two separate stochastic data augmentations $t, t' \sim T$ are applied to each example to obtain two correlated views. A base encoder network $f(\cdot)$ with a projection head $g(\cdot)$ is trained to maximize agreement in latent representations via a contrastive loss.

ea	rc	h
$\mathbf{\nabla}\mathbf{O}$	\sim	

We use random crop and color distortion for augmentation. Examples of augmentation applied to the left most images:











Google Research









f(x) is the base network that computes internal representation.

We use (unconstrained) ResNet in this work. However, it can be other networks.



ea	rc	h
$\mathbf{\nabla}\mathbf{O}$	\sim	





g(h) is a projection network that project representation to a latent space.

We use a 2-layer non-linear MLP (fully connected net).



ea	rc	h
$\mathbf{\nabla}\mathbf{O}$	\sim	



Loss function:

Google Rese

Maximize agreement using a contrastive task:

Given {x_k} where two different examples x_i and x_j are a positive pair, identify x_j in {x_k}_{k!=i} for x_i.



Original image

crop 1

crop 2

contrastive image

Let
$$sim(\boldsymbol{u}, \boldsymbol{v}) = \boldsymbol{u}^{\top} \boldsymbol{v} / \|\boldsymbol{u}\| \|\boldsymbol{v}\|$$

$$\ell_{i,j} = -\log \frac{exp(sim(\boldsymbol{z}_i, \boldsymbol{z}_j)/\tau)}{\sum_{k=1}^{2N} \mathbb{1}_{[k \neq i]} exp(sim(\boldsymbol{z}_i, \boldsymbol{z}_k)/\tau)}$$

ea	rc	h
$\mathbf{\nabla}\mathbf{O}$	\sim	

Data Augmentation for Contrastive Representation Learning



ea	rc	h
~~~	• •	

# We study a set of transformations...

Systematically study a set of augmentation



(h) Gaussian noise



(g) Cutout



#### Google Rese



(i) Gaussian blur





(j) Sobel filtering

ea	rc	h
$\mathbf{\nabla}\mathbf{O}$	$\sim$	

**Composition of augmentations are crucial** Composition of crop and color stands out!

Rotate	30.0 CTOP	22.5	20.7	4.3	9.7	6.5	2.6	13.8
	20.0	22.5	20.7	4.2	0.7	6.5	2.6	12.0
Blur	35.1	25.2	16.6	5.8	9.7	2.6	6.7	14.5
Noise	38.8	25.8	7.5	7.6	9.8	9.8	9.6	15.5
Sobel	46.2	40.6	20.9	4.0	9.3	6.2	4.2	18.8
Color	55.8	35.5	18.8	21.0	11.4	16.5	20.8	25.7
Cutout	32.2	25.6	33.9	40.0	26.5	25.2	22.4	29.4
Crop	33.1	33.9	56.3	46.0	39.9	35.0	30.2	39.2

Figure 5. Linear evaluation (ImageNet top-1 accuracy) under individual or composition of data augmentations, applied only to one branch. For all columns but the last, diagonal entries correspond to single transformation, and off-diagonals correspond to composition of two transformations (applied sequentially). The last column reflects the average over the row.

# Google Rese



(a) Without color distortion.

(b) With color distortion.

Figure 6. Histograms of pixel intensities (over all channels) for different crops of two different images (i.e. two rows). The image for the first row is from Figure 4. All axes have the same range.

ea	rc	h
$\mathbf{\nabla}\mathbf{O}$	$\sim$	

• • • • • • • • • • • • • • • •



## **Encoder and Projection Head**



ea	rc	h
~~~	• •	

Google Research A nonlinear projection head improves the representation quality of the layer before it





Google Research A nonlinear projection head improves the representation quality of the layer before it

To understand why this happens, we measure information in h and z=g(h)

What to mus dist?	Dandam guasa	Representation		
what to predict?	Random guess	h	$g(oldsymbol{h})$	
Color vs grayscale	80	99.3	97.4	
Rotation	25	67.6	25.6	
Orig. vs corrupted	50	99.5	59.6	
Orig. vs Sobel filtered	50	96.6	56.3	

Table 3. Accuracy of training additional MLPs on different representations to predict the transformation applied. Other than crop and color augmentation, we additionally and independently add rotation (one of {0°, 90°, 180°, 270°}), Gaussian noise, and Sobel filtering transformation during the pretraining for the last three rows. Both h and g(h) are of the same dimensionality, i.e. 2048.



layer when the model is asked to identify rotated variant of an image.





•••••



Loss Function and Batch Size



ea	rc	h
~~~	• •	

Norma	lized	cross	entrop

Name	Negative loss function
NT-Xent	$\left  ~~ oldsymbol{u}^T oldsymbol{v}^+ /  au - \log \sum_{oldsymbol{v} \in \{oldsymbol{v}^+, oldsymbol{v}^-\}} \exp(oldsymbol{u}^T oldsymbol{v}^+)  ight ^2$

Margin	NT-Logi.	Margin (sh)	NT-Logi.(sh)	NT-Xent
50.9	51.6	57.5	57.9	63.9

Table 4. Linear evaluation (top 1) for models trained with different loss functions. "sh" means using semi-hard negative mining.

#### Google Research

#### ov loss with adiustable temperature wor

$$\frac{\mathrm{Gradient w.r.t. } \boldsymbol{u}}{\exp(\boldsymbol{u}^T \boldsymbol{v}/\tau) \mid (1 - \frac{\exp(\boldsymbol{u}^T \boldsymbol{v}^+/\tau)}{Z(\boldsymbol{u})})/\tau \boldsymbol{v}^+ - \sum_{\boldsymbol{v}^-} \frac{\exp(\boldsymbol{u}^T \boldsymbol{v}^-/\tau)}{Z(\boldsymbol{u})}/\tau \boldsymbol{v}^-}$$



## **NT-Xent loss needs N and**

Negative loss function

 $oldsymbol{u}^Toldsymbol{v}^+/ au - \log \sum_{oldsymbol{v} \in \{oldsymbol{v}^+,oldsymbol{v}^-\}} \exp(oldsymbol{u}^Toldsymbol{v}^+)$ NT-Xent

L2 normalization with temperature scaling makes a betterloss.

Name

norm and/or temperature are changed.

$\ell_2$ norm?	au	Entropy	Contrast. task acc.	Top 1
	0.05	1.0	90.5	59.7
Vaa	0.1	4.5	87.8	64.4
res	0.5	8.2	68.2	60.7
	1	8.3	59.1	58.0
NL	10	0.5	91.7	57.2
INO	100	0.5	92.1	57.0

Table 5. Linear evaluation for models trained with different choices of  $\ell_2$  norm and temperature  $\tau$  for NT-Xent loss. The contrastive distribution is over 4096 examples.

#### Google Rese

#### Gradient w.r.t. *u*

$^{T}oldsymbol{v}/ au)$	$\Big  \; (1 - rac{\exp(oldsymbol{u}^Toldsymbol{v}^+/ au)}{Z(oldsymbol{u})})/ auoldsymbol{v}^+$ -	$-\sum_{oldsymbol{v}^-}rac{\exp(oldsymbol{u}^Toldsymbol{v}^-/ au)}{Z(oldsymbol{u})}/ auoldsymbol{v}^-$
-------------------------	----------------------------------------------------------------------------------------------------	---------------------------------------------------------------------------------------------------------

Contrastive accuracy is not correlated with linear evaluation when 12

ea	rc	h
----	----	---

#### **Contrastive learning benefits from larger batch sizes and**



Figure 9. Linear evaluation models (ResNet-50) trained with different batch size and epochs. Each bar is a single run from scratch.





### Linear evaluation 7% relative improvement over previous SOTA (cpc v2), matching fully-supervised ResNet-50.

Method	Architecture	Param.	Top 1	Top 5
Methods usir	ng ResNet-50:			
Local Agg.	ResNet-50	24	60.2	-
MoCo	ResNet-50	24	60.6	-
PIRL	ResNet-50	24	63.6	2
CPC v2	ResNet-50	24	63.8	85.3
Ours	ResNet-50	24	69.3	89.0
Methods usir	ng other architecture	es:		
Rotation	RevNet-50 $(4 \times)$	86	55.4	-
BigBiGAN	RevNet-50 $(4 \times)$	86	61.3	81.9
AMDIM	Custom-ResNet	626	68.1	-
CMC	ResNet-50 $(2 \times)$	188	68.4	88.2
MoCo	ResNet-50 $(4 \times)$	375	68.6	
CPC v2	ResNet-161 (*)	305	71.5	90.1
Ours	ResNet-50 $(2\times)$	94	74.2	92.0
Ours	ResNet-50 $(4 \times)$	375	76.5	93.2

Table 6. ImageNet accuracies of linear classifiers trained on representations learned with different self-supervised methods.





Figure 1. ImageNet top-1 accuracy of linear classifiers trained on representations learned with different self-supervised methods (pretrained on ImageNet). Our method, SimCLR, is shown in bold.





### Semi-supervised learning 10% relative improvement over previous SOTA (cpc v2), outperforms AlexNet with 100X fewer labels.

Meth

Meth Pseu VAT UD/ FixN S4L Met Instl Bigl PIR CPC Our Ours Ours

#### Google Rese

		Label fraction		
hod	Architecture	1%	10%	
		To	p 5	
hods using other label	-propagation:			
udo-label	ResNet50	51.6	82.4	
+Entropy Min.	ResNet50	47.0	83.4	
A (w. RandAug)	ResNet50	-	88.5	
Match (w. RandAug)	ResNet50	-	89.1	
(Rot+VAT+En. M.)	ResNet50 (4 $\times$ )	-	91.2	
hods using representa	tion learning only:			
Disc	ResNet50	39.2	77.4	
BiGAN	RevNet-50 $(4 \times)$	55.2	78.8	
L	ResNet-50	57.2	83.8	
2 v2	ResNet-161(*)	77.9	91.2	
s	ResNet-50	75.5	87.8	
S	ResNet-50 $(2 \times)$	83.0	91.2	
S	ResNet-50 $(4 \times)$	85.8	92.6	

Table 7. ImageNet accuracy of models trained with few labels.

ea	rc	h
$\mathbf{\nabla}\mathbf{O}$	$\sim$	

**Transfer learning** 

When fine-tuned, SimCLR significantly outperforms the supervised baseline on 5 datasets, whereas the supervised baseline is superior on only  $2^*$ . On the remaining 5 datasets, the models are statistically tied.

	Food	CIFAR10	CIFAR100	Birdsnap	SUN397	Cars	Aircraft	VOC2007	DTD	Pets	Caltech-101	Flowers
Linear evaluatio	n:											
Self-supervised	76.9	95.3	80.2	48.4	65.9	60.0	61.2	84.2	78.9	89.2	93.9	95.0
Supervised	75.2	95.7	81.2	56.4	64.9	68.8	63.8	83.8	78.7	92.3	94.1	94.2
Fine-tuned:												
Self-supervised	89.4	98.6	89.0	78.2	68.1	92.1	87.0	86.6	77.8	92.1	94.1	97.6
Supervised	88.7	98.3	88.7	77.8	67.0	91.4	88.0	86.5	78.8	93.2	94.2	98.0
Random init	88.3	96.0	81.9	77.0	53.7	91.3	84.8	69.4	64.1	82.7	72.5	92.5

Table 8. Comparison of transfer learning performance of our self-supervised approach with supervised baselines across 12 natural image classification datasets, for ResNet-50 (4×) models pretrained on ImageNet. Results not significantly worse than the best (p > 0.05, permutation test) are shown in bold. See Appendix B.6 for experimental details and results with standard ResNet-50.

* The two datasets, where the supervised ImageNet pretrained model is better, are Pets and Flowers, which share a portion of labels with ImageNet.

#### Google Research



# Conclusio

- SimCLR is a simple yet effective self-supervised learning framework, advancing state-of-the-art by a large margin.
- The superior performance of SimCLR is not due to any *single* design choice, but *a combination of* design choices.
- Our studies reveal several important factors that enable effective representation learning, which could help future research.

Code & checkpoints available in <u>github.com/google-research/simclr</u>.



ea	rc	h
$\mathbf{\nabla}\mathbf{O}$	$\sim$	

## References

- Dosovitskiy, A., Fischer, P., Springenberg, J. T., Riedmiller, M., & Brox, T. networks, arXiv preprint. arXiv preprint arXiv:1506.02753.
- preprint arXiv:1704.05310.
- by Watching Objects Move. In CVPR 2017.

Discriminative unsupervised feature learning with exemplar convolutional neural

• Bojanowski, P., & Joulin, A. Unsupervised Learning by Predicting Noise. arXiv

• Pathak, D., Girshick, R., Dollár, P., Darrell, T., & Hariharan, B. Learning Features

## References

#### **Colorization**: $\bullet$

- $\bullet$
- $\bullet$
- **Optical Flow**  $\bullet$
- and Motion Smoothness. In ECCVW, 2016.
- Zhu, Y., Lan, Z., Newsam, S., & Hauptmann, A. G. Guided optical flow learning. arXiv preprint arXiv:1702.02295.
- Others  $\bullet$
- Cruz, R. S., Fernando, B., Cherian, A., & Gould, S. DeepPermNet: Visual Permutation Learning. arXiv preprint ulletarXiv:1704.02729.
- $\bullet$ Imitation for Vision-Based Rope Manipulation. arXiv preprint arXiv:1703.02018.
- Pinto, L., Gandhi, D., Han, Y., Park, Y. L., & Gupta, A. The curious robot: Learning visual representations via physical interactions. In ECCVW 2016.

Larsson, G., Maire, M., & Shakhnarovich, G. Learning representations for automatic colorization. In ECCV 2016. Larsson, G., Maire, M., & Shakhnarovich, G. Colorization as a Proxy Task for Visual Understanding. In CVPR 2017.

• J. J. Yu, A. W. Harley, and K. G. Derpanis. Back to Basics: Unsupervised Learning of Optical Flow via Brightness Constancy

• Ren, Z., Yan, J., Ni, B., Liu, B., Yang, X., & Zha, H. Unsupervised Deep Learning for Optical Flow Estimation. In AAAI 2017

Nair, A., Chen, D., Agrawal, P., Isola, P., Abbeel, P., Malik, J., & Levine, S. Combining Self-Supervised Learning and

## **Bonus** What happens when we "super-size" self-supervised learning?

Scaling and benchmarking self-supervised visual representation learning -

Goyal, Mahajan, Gupta*, Misra* - 2019 https://arxiv.org/abs/1905.01235

**facebook** Artificial Intelligence Research

131

## Bits of information

- For ImageNet we have ~1M images and 1000 classes
  - Each image has log₂(1000) bits of information
  - Total =  $1M \times \log_2(1000)$
- For self-supervised methods we have ~1M images
  - Each image has lesser bits of information log₂(B)
  - Total =  $1M \times \log_2(B)$
  - Increase 1M to 100M?
  - Increase B to 10B?

and 1000 classes ormation

ve ~1M images tion -  $\log_2(B)$ 

132

## How large scale?

- Focus on existing **popular image-based** self-supervised methods
- Scale along three axes

#### **Problem** Complexity



facebook Artificial Intelligence Research



133

## Specific self-supervision problems



#### Jigsaw puzzles (Noorozi & Favaro, 2016)

#### **facebook** Artificial Intelligence Research



### Colorization (Zhang & Efros, 2016)

Images from the ImageNet dataset







- Use N=9 patches
- In practice, use a subset of permutations
- E.g. 100 from 9!
- Each patch is processed independently
- N-way ConvNet (shared params)
- Problem Complexity
  - Size of subset

Classify which permutation



#### 135

## Colorization



facebook Artificial Intelligence Research

### At each pixel, classify which

- Predict N=313 colors
  - Binning LAB space
  - Multi-modal loss
- Problem Complexity
  - Number of colors
  - Number of neighbors in multi-modal loss



#### 136

## How large scale?

- Scale two techniques Jigsaw and Colorization
- Scale along three axes

**Problem** Complexity "Difficulty"



AlexNet, ResNet-50



137

## Evaluating the representation





Extract "fixed" features



138

## "Investigation" task

• Train a Linear SVM on fixed feature representations

139

## "Investigation" task

- Train a Linear SVM on **fixed feature** representations
- Use the VOC07 image classification task

































plant

sheep

#### facebook Artificial Intelligence Research

sofa

train

tv

140

## Scaling on Data Axis

#### Problem Complexity "Difficulty"



**facebook** Artificial Intelligence Research Models AlexNet, ResNet-50



141

### Scaling on Data Axis



Gain for ResNet50: 10 points

Gain for AlexNet: 2 points

142

### Scaling on Data Axis



#### Gain for ResNet50: 12 points

Gain for AlexNet: 8 points

143

## Scaling on Problem Complexity

#### **Problem** Complexity "Difficulty"



facebook Artificial Intelligence Research



AlexNet, ResNet-50

**Data Size** YFCC - 100M

144
## Scaling on Problem Complexity Axis



145

## Scaling on Problem Complexity Axis Colorization – VOC07 Linear SVM



**facebook** Artificial Intelligence Research

146

# Scaling on Data and Problem

## **Problem** Complexity "Difficulty"



facebook Artificial Intelligence Research



AlexNet, ResNet-50



147



Gains along **both** data and problem axes are complementary

148

## **Our Evaluation – many tasks**



### Image classification Few-shot learning

ImageNet, Places-205, VOC'07, COCO



#### **Object detection VOC'07**

facebook Artificial Intelligence Research



#### 3D Understanding

Surface Normals – NYUv2

Navigation Gibson environment

Images from the Places, VOC07, NYUv2 and Gibson datasets





# Our Evaluation – fine-tuning vs. linear classifier



#### Fine-tune all layers

A good representation transfers with little training

**facebook** Artificial Intelligence Research



Linear classifier

150

## **Object Detection**





### Image classification Few-shot learning

ImageNet, Places-205, VOC'07, COCO

#### **Object detection VOC'07**

facebook Artificial Intelligence Research



person : 0.977







#### **3D Understanding**

Surface Normals – NYUv2

Navigation Gibson environment





# **Object Detection**

- Fast R-CNN (Girshick et al., 2015)
  - Same optimization parameters for all methods (including supervised)
  - No "bells and whistles"
  - Use VOC'07



#### **Object detection**

Image from the VOC07 dataset





## **Object Detection**

#### Initialization

ImageNet Supervised

Places Supervised

Jigsaw ImageNet-1k

Jigsaw ImageNet-22k

Jigsaw YFCC100M

facebook Artificial Intelligence Research



## VOC07 test set. Fast R-CNN ResNet50

	Train Set			
	VOC07+12	VOC07		
•	76.2	70.5		
	74.5	67.2		
within	68.3	61.4		
	75.4	69.2		
	73.3	66.6		



## **Object Detection - Training Rol heads only** VOC07 test set. Fast R-CNN ResNet50

#### Initialization

ImageNet Supervised

**Places Supervised** 

Jigsaw ImageNet-1k

Jigsaw ImageNet-22k

Jigsaw YFCC100M

facebook Artificial Intelligence Research









### Image classification Few-shot learning

ImageNet, Places-205, VOC'07, COCO



#### **Object detection** VOC'07

facebook Artificial Intelligence Research



#### **3D Understanding**

**Surface Normals – NYUv2** 

Navigation Gibson environment





- Predict surface normals on NYU-v2
  - Same optimization parameters for all methods (including supervised)
  - **PSPNet Architecture**
  - Train last few layers only (res5 onwards)



facebook Artificial Intelligence Research

### Input



#### Output

Image from the NYU dataset



#### Initialization

ImageNet Supervised

Places Supervised

Jigsaw ImageNet-1k

Jigsaw ImageNet-22k

Jigsaw YFCC100M

facebook Artificial Intelligence Research

Median Error	% correct within 11.25 ⁰
17.1	36.1
14.2	41.8
14.5	41.2
13.4	43.7
13.1	44.6



157



### Image classification Few-shot learning

ImageNet, Places-205, VOC'07, COCO



#### **Object detection VOC'07**

facebook Artificial Intelligence Research



#### **3D Understanding**

Surface Normals – NYUv2

Navigation **Gibson environment** 



158

# Visual Navigation

- Visual navigation in the Gibson environment
  - Method from Sax et al., 2018
  - Fixed ConvNet features



Gibson Environment

159

## Visual Navigation



## Few shot learning



### **Few-shot learning**

Places-205, VOC'07



#### **Object detection** VOC'07

facebook Artificial Intelligence Research



#### **3D Understanding**

Surface Normals – NYUv2

Navigation Gibson environment





# Few shot learning

- k-shot learning
  - Use VOC'07/Places205 classification
  - K labeled examples per class
  - Train linear SVMs



### Few-shot learning

Image from the Places dataset



### Few shot learning **VOC07** 80 mageNet-1k Supervised 60 Jigsaw ImageNet-22k mAP 40 Jigsaw YFCC-100M 20 Random 01 32 2 16 8 4 Num. Labeled samples

#### **facebook** Artificial Intelligence Research

### Self-supervised representations are not as sample efficient





## Few shot learning - VOC07







-supervised conv4 -yfcc conv4

VOC07 low-shot conv3



#### VOC07 low-shot conv5

-supervised conv5 -yfcc conv5

## Few shot learning - Places205



places205 low-shot svm conv4



places205 low-shot svm conv5 50. 37.5 25. 12.5 0. k=1 k=2 k=4 k=8 k=16 k=32 k=64 k=96

—supervised conv5 —yfcc conv5

#### 165

## Image Classification



### Image classification

Places-205, VOC'07



#### **Object detection VOC'07**

facebook Artificial Intelligence Research



#### **3D Understanding**

Surface Normals – NYUv2

Navigation Gibson environment





## Linear SVMs on VOC07

VOC2007 SVM classification. ResNet50						
Init	conv1	stage1	stage2	stage3	stage4	
ImageNet supervised	24.49	47.75	60.54	80.36	87.95	
ImageNet jigsaw	27.09	45.73	56.61	64.51	57.17	
Imagenet14M jigsaw	23.46	46.72	58.52	71.76	64.92	
YFCC100M jigsaw	18.98	46.71	57.77	71.21	64.34	

facebook Artificial Intelligence Research

Deeper self supervised layers are less transferable Hypothesis – Problem is not "complex" enough.

167

# SGD-based Linear Classifiers on Places205

Places205 linear classification. ResNet50

Init	conv1	stage1	stage2	stage3	stage4
ImageNet supervised	14.84	32.59	42.06	50.83	52.49
ImageNet jigsaw	15.079	28.753	36.825	41.232	34.364
Imagenet14M jigsaw	13.973	29.462	36.656	41.721	36.254
YFCC100M jigsaw	10.252	29.672	39.565	44.866	38.219

- YFCC is a good mixture of both.

facebook Artificial Intelligence Research • Gap between ImageNet and self-supervised methods is smaller.

Places has scenes while ImageNet is object centric

168

## **Evaluation: Main Lessons**

- Evaluation on multiple tasks is essential
- Evaluation with fixed features or at least same hyper-parameters
- Evaluate sample efficiency of representations

## Itial least same hyper-parameters sentations

# What's missing from self-supervised methods?

- Complex problems, big data and deeper models
- Current self-supervised methods do not seem to learn high level representations
- Sample efficiency



## Image classification Few-shot learning

ImageNet, Places-205, VOC'07, COCO

#### facebook Artificial Intelligence Research



#### Object detection VOC'07

# Thanks!



## 3D Understanding

Surface Normals – NYUv2

Navigation Gibson environment

